# EXAMPLE

- What are some good predictors of whether you will stay in or go out for this weekend?

- Factors : Weather, your spouse in town, important game on TV, do you have something urgent to do at work.

- All these factors usually have some order of importance.

# WHAT ARE DECISION TREES?

- In Machine Learning decision tree predicts the outcome given the values of the input variables.

- It's a Supervised learning algorithm.

- Decision Tree helps in making decision in structured manner and visually represents all the variables involved and the consequences in each case.

# Advantages and Disadvantages

1. Easy to interpret and perfect for visual representation.

2. It works both with categorical and numerical features.

3. Nonlinear parameters don't affect its performance.

4. Little effort is required in data preparations.

1. Overfitting as algorithm tends to pick up the noise in the data.

2. High Variance as model tends to get unstable with little variation in data.

3. Low bias which makes the model difficult to work with new data.

# WORKFLOW OF DECISION TREES

**Keywords of Decision Tree**

- The topmost feature would be called as **root node**

- Then we have other features as **branches**

- The endpoints are called **leaf nodes.**

**Entropy :** The measure of randomness and unpredictability in the dataset.

Range(0 to 1)

$$E = -\sum_{i}^{C} p_i \log_2 p_i$$

**Information Gain** is measure of entropy after the dataset has been split.

Range(0 to 1)

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

**Gini Impurity** is used for calculation of purity of a split.

Range (0 to 1)

$$G = \sum_{i=1}^{C} p(i) * (1 - p(i))$$

# How the Decision Tree works
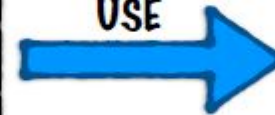
# DECISION TREE LEARNING ALGORITHMS BASED ON RECURSIVE PARTITIONING

EACH HAS A SLIGHTLY DIFFERENT WAY OF ARRIVING AT THE BEST ATTRIBUTE (OR) MEASURING THE HOMOGENEITY OF A SUBSET

ID3
C4.5 — USE → INFORMATION GAIN

CART — USES → GINI IMPURITY

CHAID — USES → STATISTICAL SIGNIFICANCE

# RANDOM FOREST ALGORITHM

# Bagging (Ensemble Technique)

- Also known as Bootstrap Aggregation.
- Row Sampling with Replacement is performed (Bootstrap).
- Each Model trains and predicts an outcome and the majority is taken as final output (Aggregation).

# Random Forest

- Random forest, like it names implies consists of large number of individual decision trees that operate as an ensemble.

- Each individual tree has spits out a class prediction and class with most votes becomes our final model prediction.

- In Random Forest row sampling and feature sample is used which makes its a powerful algorithm.

# Application of Random Forest

**Remote Sensing**
Used in ETM devices to acquire images of the earth's surface.

Accuracy is higher and training time is less

**Object Detection**
Multiclass object detection is done using Random Forest algorithms

Provides better detection in complicated environments

**Kinect**
Random Forest is used in a game console called Kinect

Tracks body movements and recreates it in the game

# Why Random Forest?

**No overfitting**

Use of multiple trees reduce the risk of overfitting

Training time is less

**High accuracy**

Runs efficiently on large database

For large data, it produces highly accurate predictions

**Estimates missing data**

Random Forest can maintain accuracy when a large proportion of data is missing

**Boosting** (Ensemble technique)

- ADA Boost (Adaptative Boost)

- Gradient Boost

- XGB  Boost

# Performance Metrics

# ROC and AUC Curve

ROC : Receiver Operating Characteristic is mostly used to visualize binary classifier.

AUC : Area Under Curve is measure of ability of a classifier to distinguish and be the summary of ROC curve .

The ROC is plotted using the True Positive Rate and the False Positive Rate.

# ROC and AUC

- Its plotted is given as Sensitivity vs 1- Specificity.

- Sensitivity is True Positive Rate which is ratio of True positives and sum of True positives and False negatives.

- Specificity is False Positive Rate which is of False Positives and sum of True negatives and False positives.

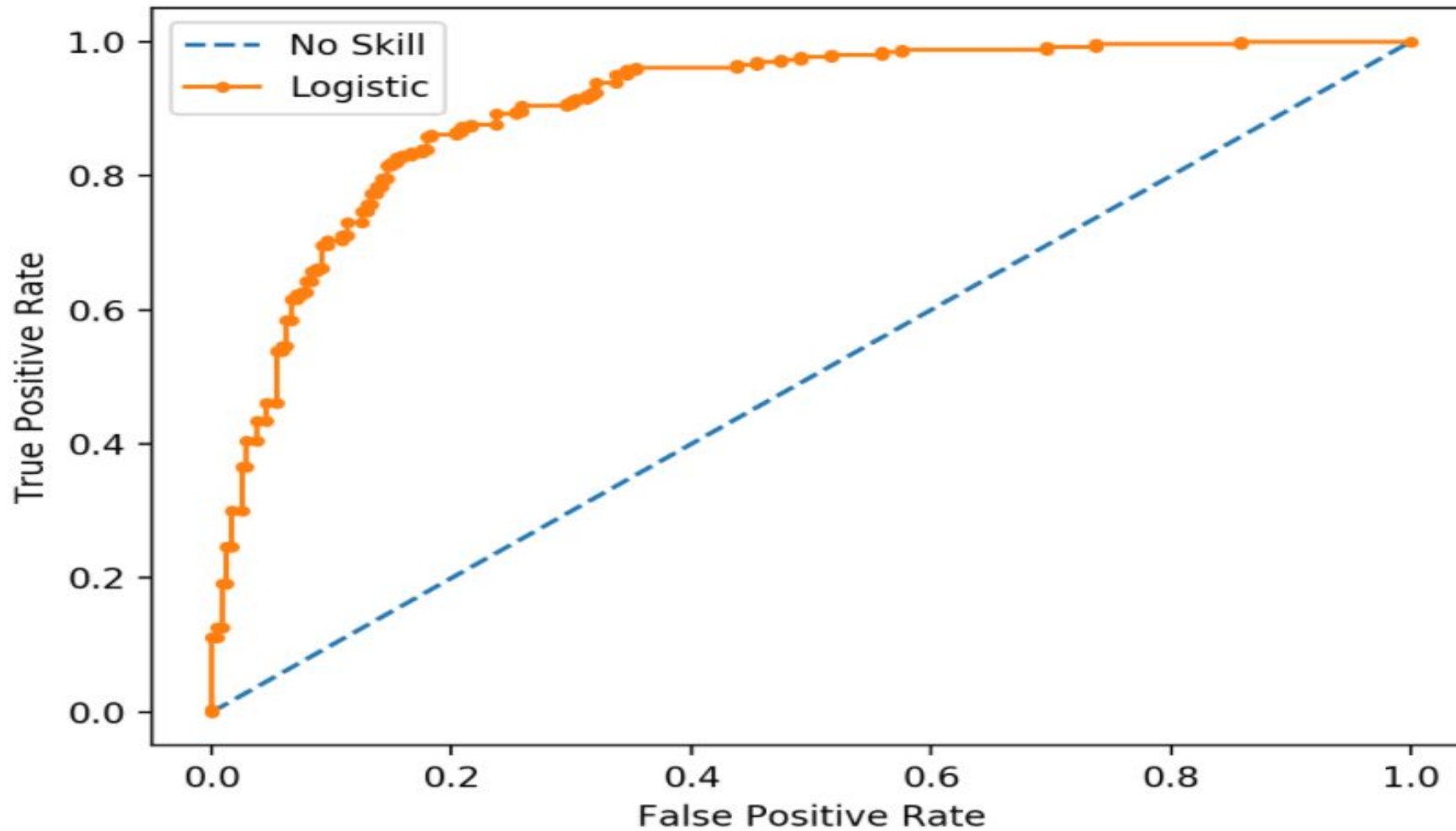ROC Curve Plot for a No Skill Classifier and a Logistic Regression Model

# Industry Use cases of Decision Trees

- Commonly used in data mining .

- Loan Approval

- In Finance sector, forecasting future outcomes and assigning probabilities to those outcomes

# Industry Use cases of Random Forest

- Banking Sector : Banking sectors consists of most users. Used for Fraud Detection

- Medicines needs complex combination of chemicals. Random forest comnes handy in this scenario.

- Stock Market : Stock behaviour analysis can be done using Random Forest.

# Summary

- Decision Tree model is like a white box.

- Decision tree performs well in case of non linear data.

- Decision tree tends to overfit.

- Overfitting can be mitigated using various Ensemble techniques

- Random Forest is an ensemble of decision trees.

- Random forest widely popular because of its missing value handling and no overfitting.

THANK YOU
Hope You Got To Learn Something!!!